

UNIVERSITY OF TARTU
FACULTY OF MATHEMATICS AND COMPUTER SCIENCE
Institute of Computer Science

Stenver Jerkku

Parallel Wilcoxon signed-rank tests

Bachelor's Thesis (6 ECTS)

Supervisor: Sven Laur, PhD

Tartu 2014

Paralleelsed Wilcoxon astaku testid

Lühikokkuvõte: Statistilisi teste kasutatakse tuvastamiseks kuidas erinevad eksperimentaalsed stiimulid mõjutavad uuritavaid suurusid. Antud töös uuritav Wilcoxon test on üks vähestest statistilistest testidest, mida saab kasutada juhul kui grupi sees olev loomulik varieeruvus pole normaaljaotusega. Selliste testide kasutamine on tavapärane bioloogiliste andmete uurimisel. Sellest tulenevalt kasutatakse seda testi Bioinformaatika, Algoritmika ja Andmekaeve grupi poolt geenide uurimiseks, analüüsimiseks ja andmekaeveks, bioloogiliseks andmekaeveks ja muudeks ülesanneteks.

Praegused implementatsioonid Wilcoxon astaku testist on optimeerimata ja aeglased. See projekt vaatab Wilcoxon testi põhiomadusi ning uurib kuidas selle implementatsiooni optimeerida. Selleks, et implementatsiooni teha täpsemaks, uuritakse Wilcoxon statistiku ja Gaussi jaotuse seost. Selleks, et implementatsiooni kiiremaks teha, kasutatakse dünaamilise programmeerimise meetodeid, et säästa arvutusaega. Optimeerimisega tehti teste nii kiiremaks kui ka täpsemaks.

Antud töös loodi täpne ja kiire Wilcoxon testi implementatsioon C++ jagatud teek. Selle projekti skoobis on ka nimetatud teegi integreerimine käsureaga ja GNU-R projektiga. Tänu enda jagatud teegi olemusele, on seda lihtne kasutada ja implementeerida ka teistes tööriistades.

Võttesõnad: Wilcoxon astaku test, paralleelsed testid, statistiline testimine, numbriline approksimatsioon, GNU-R, C++.

Parallel Wilcoxon signed-rank tests

Abstract: Statistical tests are used to find out if some sort of experimental stimulation affects observable features. In this paper we researched Wilcoxon signed-rank test which is one of the few statistical tests that can be used when the natural variation inside the group is not normally distributed. The test is used by Bioinformatics, Algorithmics and Data mining group research for gene regulation, gene expression data analysis, biological data mining and others.

The current implementations of the Wilcoxon signed-rank tests are slow and unoptimized. This project looked into the foundations of Wilcoxon signed-rank test, its current implementations and how to optimize it. In order to make the implementation more accurate, the relationship between Wilcoxon statistic and Gaussian approximate was observed. In order to make the implementation faster, some dynamic programming methods were used to save computation time. The purpose of optimizing was to make it more accurate and speed up the test running.

In this project an accurate and fast Wilcoxon test shared library was created. In the scope of this project, the library was integrated with two tools - command line and GNU-R. Due to the nature of shared library, it will be easy integrate the library with any other tools one might desire.

Keywords:Wilcoxon signed-rank test, parallel tests, statistical testing, numerical approximation, GNU-R, C++.

Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 7 |
| 1.1 | NetCDF | 7 |
| 1.2 | Current BIIT procedures | 8 |
| 1.3 | List of thesis contributions | 8 |
| 2 | Wilcoxon signed-rank test | 10 |
| 2.1 | Statistical Tests | 10 |
| 2.1.1 | Test Statistic | 10 |
| 2.2 | Hypothesis | 11 |
| 2.2.1 | Null hypothesis | 11 |
| 2.2.2 | Alternative hypothesis | 11 |
| 2.2.3 | P-value | 12 |
| 2.3 | Wilcoxon test | 13 |
| 2.4 | Wilcoxon test assumptions | 13 |
| 2.5 | How to compute | 14 |
| 2.6 | Example | 15 |
| 2.7 | P-value assignment | 15 |
| 2.7.1 | Exact computation of P-value | 17 |
| 2.7.2 | The V and P table formulas | 18 |
| 2.7.3 | The Gaussian approximation of P table | 19 |
| 2.7.4 | Code sample | 19 |
| 2.7.5 | Bonferroni correction | 20 |
| 3 | Approximating the p-value | 21 |
| 3.1 | Motivation | 21 |
| 3.2 | When can we approximate | 22 |
| 3.3 | The behavior of relative error | 28 |
| 3.4 | Finding out when can we approximate p-value | 31 |
| 4 | Further optimizations | 32 |
| 4.0.1 | dose-response curve | 33 |
| 4.1 | Linear approximation | 34 |
| 4.1.1 | Goal | 34 |
| 4.1.2 | Linear approximation | 34 |
| 4.1.3 | Linear interpolation | 35 |
| 4.1.4 | Algorithm for dropping intermediate points | 36 |
| 4.2 | Optimization summary | 37 |
| 4.3 | Further optimizations | 38 |

| | | |
|----------|----------------------------------|-----------|
| 5 | The implementation | 39 |
| 5.1 | Library description | 39 |
| 5.2 | Repository description | 39 |
| 6 | Conclusion | 41 |

1 Introduction

The biologists of BIIT research group are conducting a variety of experiments. They try to measure, compare and test different attributes of a living organism. These attributes might be, but not limited to blood pressure, protein amount in blood, amount of RNA, purity of urine, brainwaves etc. These attributes can be repeatedly measured on the same subject and the results are never exactly the same. There is always a little noise and variance between the samples. Biologists often experiment on the control group and try to manipulate with these observed attributes. They take samples before and after stimulating the observable attribute with some stimulant. This is called the case-control study.

However, even without external stimulation the two results are never the same - human body constantly has minor fluctuations on its observable features. For example, if you measure white blood cell count two times in a row, you almost never get the exact same result. Also, different experiments can have wildly different assumptions, so you cannot simply use a constant error margin on all your tests. For example, the blood pressure measurements before and after jogging can probably have a lot bigger changes than brainwave measurements. Therefore, the biologists need to use statistical tests to find out which case-control studies are actually significant and which are simply random fluctuations.

There are many case-control tests available. For example paired Students t-test, t-test for matched pairs and Wilcoxon signed-rank test. The Wilcoxon signed-rank test is used when the measurement values are not normally distributed and when the experimental conditions are fixed. In practice, this means that the histogram of measurements is either asymmetrical or it does not resemble to the bell shape.

It is common to assume that properly scaled gene expression measurements follow Gaussian distribution, while quantities of proteins and metabolites are assumed to have non-Gaussian distribution.

1.1 NetCDF

The BIIT group holds statistical data in NetCDF file format. NetCDF is an open source standard set of data formats, programming interfaces, and software libraries that help read and write scientific data files [RGSH11].

Data can be held in structured manner in a NetCDF file. You can define dimensions, name them, put variables in the dimensions and later retrieve or change them. This is useful, for example, to hold matrix-like data in a file and have fast lookups. In addition, you can define helper dimensions for the matrix in case of needing the extra information. NetCDF file format is used by the BIIT research group to hold gene data and will be given as an input file for the program.

1.2 Current BIIT procedures

Currently the BIIT group has 3 ways of analyzing its data with Wilcoxon test. However, all of these solutions ultimately use GNU-R implementation.

GNU-R. They could use GNU-R project to make statistical analyses. To achieve this purpose, they manually read the NetCDF file to R and then call the built in R function `wilcoxon.test` to use Wilcoxon test. The problem with R Wilcoxon test is that it is slow. If you run thousands of tests in a row, the results might come after few minutes, or even hours, depending on the processing power of the computer. One of the goals of this paper and project is to make Wilcoxon test run on command line in mere seconds.

Web interface. They could use a web interface, which can take in certain arguments and a command line script. Subsequently it will use whatever command line tools it has available and visualizes the output. The efficiency problems in servers are even more important, since inefficient solutions waste precious server time.

Command line. They could also use the command line directly which has a variety of tools installed. These are needed to create large data processing pipelines. The Wilcoxon test will ultimately call R Wilcoxon function.

1.3 List of thesis contributions

I would like to thank Sven Laur who was a tremendous help with the theoretical framework, analysis and refining this paper.

Analyzing the Wilcoxon signed-rank

I described how the statistical test works and also gave an in-depth overview of the fundamentals of the Wilcoxon test.

Analyzing the approximation of p-value

I questioned the common assumption, which is made while using Wilcoxon test - meaning when the sample size is at least ten, then you can already approximate the p-value. The research revealed some interesting results, which proved that when running a single test, this sample size is adequate. On the other hand, when running parallel tests, the required sample size will get larger, as the number of tests grows.

Optimizing the Wilcoxon test

I started with the idea of hardcoding the accurate p-value table. After analyzing the results of Chapter 3, I further optimized the solution, which involved experiments on multiple algorithms for the most efficient solution and approximating the accurate p-value table itself through more accurate algorithms than Gaussian distribution.

Implementation

I implemented the Wilcoxon test algorithm which involved all the optimization tricks learned during the research phase of this paper. The implementation itself is an open sources C++ shared library, which is easily extendable to other open source projects.

Integration

I integrated the shared library with GNU-R and command line. The integrations are split into separate sub-projects in the thesis repository to enforce decoupling and good extendibility of the original shared library.

2 Wilcoxon signed-rank test

In this section, we discuss a particular statistical test called Wilcoxon signed-rank test. Wilcoxon test can be used to make sure that something has an effect on the measured samples. For example, biologists could measure three patients' white blood cell levels, create some external stimulation on the patients and measure the white blood cell levels again. The Wilcoxon signed-rank test shows whether the external stimulation had any statistical effect on the white blood cell levels.

2.1 Statistical Tests

If the initial research has raised a hypothesis about the data then it is desirable to know if the hypothesis is true or not. To find out if the hypothesis is true, statistical tests are used.

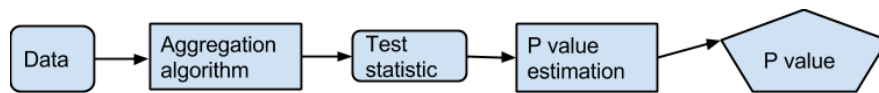


Figure 1: Graphical representation of statistical test flow

Statistical tests are usually built by first gathering the necessary raw data. For example, make 80 persons run for 30 minutes and measure blood pressure before and after running. Then the data should be organized by some kind of aggregation algorithm which also might filter out samples that are not related to the hypothesis in question. For example, sorting and grouping the blood pressure results by gender and weight, in ascending order. Then filter out any group that had too few samples. A test statistic is chosen and calculated afterwards. After that a p-value can be estimated, which tells if the observable feature is interesting or not. For example, if running plays a role in raising a person's blood pressure.

2.1.1 Test Statistic

The test statistic is a scalar function of all the observations, which summarizes the data by a single number. This value is used to estimate a p-value and accept or reject the hypothesis. Commonly, the test statistic is either one-sample, two-sample or paired statistic, depending on the test [RFDE14].

A one-sample test investigates whether the data collectively satisfies some hypothesis, e.g. the data is generated by Gaussian distribution or has specific theoretical mean value(expected value).

A two-sample test compares two sub-populations of data, e.g. persons with high blood pressures and persons with low blood pressures.

A paired test characterizes the change of before and after measurements, e.g. does some novel medical drug treatment work. The differences between measured observable case and control experiment pairs comes from the same distribution. For instance, we might use paired statistical distance to estimate whether the differences between paired measurements are symmetrical or not. If the differences are symmetrical, then a positive difference x is equally likely as a negative difference $-x$.

2.2 Hypothesis

2.2.1 Null hypothesis

In statistics, a null and alternative hypothesis are raised to find out if the test bears any interesting results. A null hypothesis in general terms means that the test results are boring or that there is no effect. It means there is no relationship between the two observed features. For example, when measuring blood pressure before and after running, the null hypothesis would state that blood pressure will not change when running [MTrio01].

The null hypothesis also defines the data distribution. For example, it might define that data is generated by a Gaussian distribution. Since null hypothesis is the boring result, i.e., that's what the data set usually looks like, then generally the data will be distributed according to the null hypothesis. Since the null hypothesis defines the data distribution, it will also define the test statistic distribution. Finally, since the null hypothesis defines the data distribution and test statistic distribution, it will also play a role in estimating the p-value for the test statistic.

For example, if the observations x_1, x_2, \dots, x_n are generated by coin flip, then the test statistic $T = x_1 + x_2 + \dots + x_n$ has a binomial distribution. With n number of trials, one side of the coin always has a probability of $\frac{1}{2}$ to occur with every flip. If test statistic is defined as $T = x_1(1 - x_1) + x_2(1 - x_2) + \dots + x_n(1 - x_n)$, then it is constantly zero and thus does not hold any relevance. The exact form of the test statistic determines the properties of the test.

2.2.2 Alternative hypothesis

The alternative hypothesis is the opposite of to the null hypothesis. In simple terms, one could think of it as true or interesting result. It means that the two measured samples have a relationship between them, either in direct or indirect ways. A good example to explain these two is the court verdict. If null hypothesis is kept, then the suspect has not been found guilty. There was not enough evidence. On the contrary, if null hypothesis was rejected, the suspect was found guilty, because there was enough evidence to make that decision. Although the alternative

hypothesis could be anything that is not a null hypothesis, some alternatives to the null hypothesis are easier to distinguish from the null hypothesis.

When the null hypothesis is true, then the test statistic has a well-defined probability distribution. This means that one can predict pretty well the possible outcome of subset of samples. For example, the alternative hypothesis can be separated from the null hypothesis if the distribution of the test statistic is vastly different from the null hypothesis. The alternative hypothesis corresponding to the null hypothesis causes the test statistic to be at the edges of a probability distributions, which makes them hard to predict.

2.2.3 P-value

P-value is estimated from the test statistic. It is used to find out if the test data is interesting. A certain threshold, called significance level (usually 0.05), is chosen. If the p-value is within significance level, then null hypothesis is rejected. If the p-value is below significance level, then null hypothesis is kept and the data falls into distribution as defined by null hypothesis. A variety of algorithms have been thought for p-value estimation, depending on the type of data, goals of the tests and even the amount of data. Statistics need to make the right choice between the algorithms by taking all these assumptions into account.

The p-value can be either one-sided or two-sided. Given a test statistic T_0 which we computed from observed data. A one-sided p-value shows the probability that we can get a value T that is larger or equal than the T_0 in the distribution defined by H_0 . This can be seen visually in Figure 2.

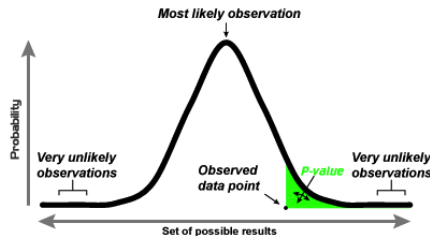


Figure 2: One sided p-value representation. A p-value(shaded green area) is the probability of an observed (or more extreme) result arising by chance [ReWi14].

A two-sided p-value shows the probability that we can get a value T that is larger or smaller than the T_0 in the distribution defined by H_0 . This can be seen visually in the Figure 3.

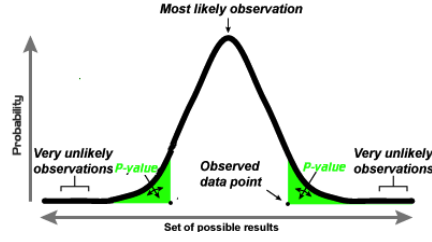


Figure 3: Two sided p-value representation. A p-value(shaded green area) is the probability of an observed (or more extreme) result arising by chance [ReWi14].

2.3 Wilcoxon test

The Wilcoxon signed-rank test is a statistical hypothesis test that is used when comparing two related samples, matched samples, or repeated measurements on a single sample. This means that with Wilcoxon test you can look at some measurable feature, e.g. blood pressure. If two or more experiments in different conditions have been made on it, e.g. gene experiments and their confirmation results, Wilcoxon test can be used to see if the changes of the measurement levels are relevant. The null hypothesis H_0 means that the difference between pairs is zero. The alternative H_1 means that it is different from zero [FWil45].

It was popularized by Sidney Siegel in his book “Nonparametric statistics - for the behavioral sciences”. Sidney used the symbol T , to denote test statistic [SSieg56]. Because of this, the test is sometimes referred to as the Wilcoxon T test. However, in this paper the more common W is used to represent the Wilcoxon test statistic value.

2.4 Wilcoxon test assumptions

1. Two measurements can always be compared and thus ordered. For example, measurements are real values or ordered categorical values.
2. Measurements can be naturally paired into case-control or before-after experiments. For example, measurements of patient’s white blood cell levels before and after treatment.
3. All measurement pairs are independent from other pairs. For example, measurement of one individual does not influence the other. In medicine studies, the patients are sampled randomly. There is no evident procedure in selecting patients or no planned action to organize a study so one could get the “desired” result.

2.5 How to compute

Let N be the sample size, the number of pairs. Then we can use the following variable to compute the W test.

For $i = 1, \dots, N$, let $x_{1,i}$ and $x_{2,i}$ of the same quantity in case and control group denote the measurements.

Let W be the Wilcoxon test statistic.

Let z -score be the Wilcoxon test standard score.

Let N_0 be the number of pairs from which onward we can approximate the p-value. This value is currently undefined and is the focus of this paper. It will be discussed in the future chapters.

The computation of the statistic W is organized into five steps. The last step diverges into two possible paths and the choice of the path that is chosen depends whether $N \geq N_0$.

1. For $i = 1, \dots, N$, calculate $|x_{2,i} - x_{1,i}|$ and $\text{sign}(x_{2,i} - x_{1,i})$, where $\text{sign}(x)$ is the sign function.
2. Order the N pairs from smallest absolute difference to largest absolute difference, $|x_{2,i} - x_{1,i}|$.
3. Rank the pairs, starting with the smallest as 1. Ties receive the rank equal to the average of the ranks they span. Let R_i denote the rank.
4. Calculate the test statistic W , the absolute value of the sum of the signed ranks.

$$W = \left| \sum_{i=1}^N [\text{sign}(x_{2,i} - x_{1,i}) R_i] \right| \quad (1)$$

5. As N increases, the sampling distribution of W converges to a Gaussian distribution. Where N_0 depends on how accurate you want your results to be.

(a) For $N \geq N_0$, a z -score can be calculated as

$$z = \frac{W - 0.5}{\sigma_w} \quad (2)$$

$$\sigma_w = \sqrt{\frac{N(N+1)(2N+1)}{6}} \quad (3)$$

If $Z > Z_{critical}$ then reject H_0 , where $Z_{critical}$ is calculated with Gaussian distribution.

(b) For $N < N_0$, W and N is used to calculate accurate P value as described in Section 2.7.2.

If $P \geq P_{critical}$, N then reject H_0

2.6 Example

Given pairs of measurements $(6, 8), (2, -3), (-3, 3), (1, 3)$.

1. Calculate absolute values and signs of the pairs. For example, $(6, 8)$ value becomes 2 and sign 1.
2. Order the pairs. Ties receive the rank equal to average of the ranks they span. Let R_i be the rank value.
3. Calculate the test statistic $W = 1, 5 - 3 + 4 + 1.5 = 4$.
4. Since N_r is very small, use a table to look up the p value. The calculation of P table is shown in Section 2.7.2.
 $P(4, 4) = 0.3125$
 Since $0.3125 > 0.05$, reject H_1

The calculations can also be seen on Table 1.

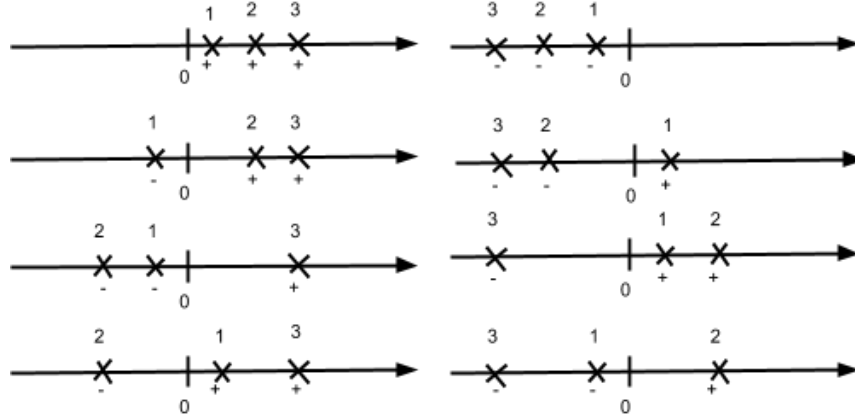
| i | $x_{1,i}$ | $x_{2,i}$ | sign | abs | R_i | sign * R_i |
|---|-----------|-----------|------|-----|-------|--------------|
| 1 | 6 | 8 | 1 | 2 | 1.5 | 1.5 |
| 2 | 2 | -3 | -1 | 5 | 3 | -3 |
| 3 | -3 | 3 | 1 | 6 | 4 | 4 |
| 4 | 1 | 3 | 1 | 2 | 1.5 | 1.5 |

Table 1: Initial data

2.7 P-value assignment

Let us consider Wilcoxon test if we have 6 measurements. Then there are three differences as shown in in Figure 4. As the test considers only the sign and order of the measurements, then only 8 configurations are possible, as shown in Table 2. Due to the symmetry of the possible configurations under the distribution of the null hypothesis, all 8 configurations can be proven equally probable.

Now if the W value in our data is +2, then there are 3 configurations that have larger or equal W value. This can be seen from the Table 2. Then the one-sided



$$W = 1 + 2 + 3 = 6$$

$$W = -1 + 2 + 3 = 4$$

$$W = -1 - 2 + 3 = 0$$

$$W = 1 - 2 + 3 = 2$$

$$W = -1 - 2 - 3 = -6$$

$$W = 1 - 2 - 3 = -4$$

$$W = 1 + 2 - 3 = 0$$

$$W = -1 + 2 - 3 = -2$$

Table 2: The figure shows all possible configurations of differences. The equations show values of the test statistic corresponding to the configurations.

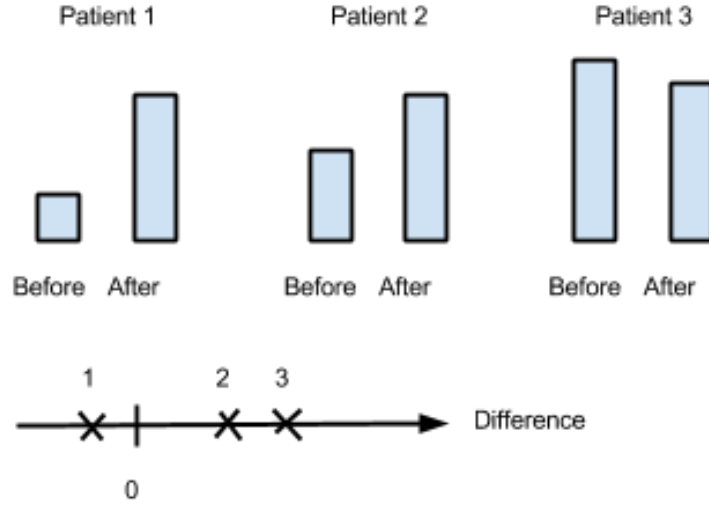


Figure 4: Patient example

p-value is $\frac{3}{8}$. The two-sided p-value is $\frac{3}{4}$. The same procedure is needed to assign p-value.

Wilcoxon test neglects the actual value of the tests. To estimate the p-value, we only need to know the W value and number of tests N . W statistic will always create the same configuration for each possible N .

2.7.1 Exact computation of P-value

In general, we get 2^n equally probable configurations of differences and assigning a p -value can be reduced to the following problem. Given N ranks, we assign $+$ or $-$ sign to each rank randomly and each sign has equal probability to be assigned. We can consider all 2^n possible sign assignments and find out in how many ways we can assign the ranks to get a certain sum k . If we do this for every possible k , then a V -array is formed. It signifies all the possible random signed rank combination sums in an N sized ranked array, where k shows how many times a certain sum was achieved. If we calculate V -array for every possible number of ranks up to the fixed N , then a V -table is formed which we can use through $V_{N,k}$ to find out how many times a certain sum k was achieved.

From this table, the probability that $W \geq k$ can be calculated. This value will show that given N ranks with random signs, how probable it is that a sum k is the result. If this process of calculating values is done repeatedly for every possible N

and k , a P -table forms. The P -table shows the probability of every possible sum k appearing given N ranks with random signs. Let $P_{N,k}$ be the function to get the desired P value from that table.

If the P -table gets large enough, it will start taking the shape of Gaussian distribution. Given enough ranks, a Gaussian distribution can be used to approximate the P -table.

2.7.2 The V and P table formulas

Given N which shows us the number of ranks starting from 1 and k which shows us the sum that we are interested in. Recursively apply this algorithm until you reach to the $N = 1$.

$$V_{N+1,k} = V_{N,k-1} + V_{N,k+N+1} \quad (4)$$

There are some additional conditions to the formula:

$$V_{N,k} = 0 \quad \text{if} \quad k < \frac{-N(N+1)}{2} \quad \text{or} \quad k > \frac{N(N+1)}{2} \quad (5)$$

$$V_{N,k} = 1 \quad \text{if} \quad k = \frac{-N(N+1)}{2} \quad \text{or} \quad k = \frac{N(N+1)}{2} \quad (6)$$

With these formulas, the probability that W is a certain value can be calculated. Given V table, N and k , we can calculate the P as follows.

$$P(T) = \frac{1}{2^n} \cdot \left| \sum_{K=T}^{\infty} W_{N,K} \right| \quad (7)$$

To get the P table, we go through all N and K values that we are interested in.

As a concrete example, consider the case where $N = 2$, then ranks are 0, 1, 2. One can combine them as follows.

$$\begin{aligned} &0 + 1 + 2 \\ &0 + 1 - 2 \\ &0 - 1 + 2 \\ &0 - 1 - 2 \end{aligned}$$

This gives us one way to get 3, one way to get 1, one way to get -1 and one way to get -3 . It can now be calculated that there is 0.25% probability that our value is 3 or lower, and 0.5% probability that our value is 1 or lower.

2.7.3 The Gaussian approximation of P table

As mentioned above, if the number of ranks N gets large enough, then Gaussian distribution can be used. To use that, the cumulative density function $F(x)$ is usually used. A simple Gaussian function $F(x, 1)$, where x is the statistic z -value, e.g. standard score and 1 is the sigma, will give approximately the same results as the accurate P table.

2.7.4 Code sample

A python code to calculate the V table is as follows.

```
def V(n, k, vTable):
    if((n == 1 and (k == 1 or k == -1)) or (n == 0 and k == 0)):
        :
        return 1
    if((k < -(n * (n + 1) / 2)) or (k > (n * (n + 1) / 2))):
        return 0
    n = n - 1

    kIndex = k + kSize

    leftValue = 0 if kIndex - n - 1 < 0 or kIndex - n - 1 >=
        len(vTable[n]) else vTable[n][kIndex - n - 1]
    rightValue = 0 if kIndex + n + 1 < 0 or kIndex + n + 1 >=
        len(vTable[n]) else vTable[n][kIndex + n + 1]

    return leftValue + rightValue

def calculateVValues():
    vTable = []
    for n in range (nSize):
        tableRow = []
        for T0 in range (kSize * 2 + 1):
            tableRow.append(V(n, T0 - kSize, vTable))
        vTable.append(tableRow)
    return vTable
```

The python code to calculate the P table is as follows:

```
def calculatePValues(vTable):
    pTable = []
    for n in range(nSize):
        tableRow = []
        sumValue = 0
        powerValue = 2**n
        for T0 in range(kSize * 2, kSize-1, -1):
            sumValue += vTable[n][T0]
            p = sumValue / powerValue
            tableRow.insert(0, p)
        pTable.append(tableRow)
    return pTable
```

You can find a python implementation of a code sample in repository [JSpac14]

2.7.5 Bonferroni correction

If you make multiple hypothesis on a test, then you increase the risk in which you reject null hypotheses when it's actually true. The Bonferroni test helps to counteract this in a simple and naive way. For each hypothesis you make on a test, you should use a significance level M times lower than before. This ensures that the increased risk in which we can reject null hypothesis will not raise, no matter the number of hypothesis on a test. So for example, if you make M hypothesis and want a significance level α , then you should run each test at a significance level of $\frac{\alpha}{M}$. If you want to use Bonferroni correction with Wilcoxon signed-rank test, then you need to keep in mind that the approximated P value significance level needs to be much more accurate, since the Bonferroni test divides it by the amount of tests made.

3 Approximating the p-value

3.1 Motivation

The current implementations of Wilcoxon signed-rank test usually assume that p -value distribution is sufficiently close to Gaussian distribution, so that we can use Z to calculate the P value. They draw this conclusion from the assumption that the number of samples N is sufficiently large to use Z . When $N = 1$, then approximation is quite inaccurate. As number of samples increases, then the approximation becomes more and more accurate on absolute scale. When the $N \rightarrow \infty$, then the P_{approx} approximate region will almost equal the P accurate region, which also means that the region will be almost perfect Gaussian distribution. This means that when $N \rightarrow \infty$, then there is no reason not to use Gaussian distribution, as it will be almost completely accurate, except at the very edges. This allows them to use Gaussian distributions to approximate the p-value and the accurate P table calculation is often left unoptimized.

In the BIIT research group, and many others, the N is usually not sufficiently large. This means that Gaussian distribution should not be used. Also, since BIIT uses multiple hypothesis testing, the approximation must be good even at the very edges of distribution. As mentioned above, the Bonferroni correction will lower the significance level significantly as N rises. Furthermore, it is not known exactly where the N limit is, so an arbitrary number is used for that purpose. The textbooks [RLow11] currently suggest a sufficient N is 10. In this work we studied the question whether this assumption is justified in detail.

The approximation is computationally cheap and the preferred method to be used over accurate P table, which is computationally expensive. Because of this, one focus of the paper is finding out how good the approximation of Gaussian distribution is and finding out the exact N where we can use approximation instead on accurate P table while maintaining a high enough accuracy. Let N_0 be the number of samples that allows us to use approximation accurately.

The bottleneck of the other implementations of the Wilcoxon test is, when N is low, but the number of tests is high. In this case they calculate the P and V tables every time the test is run, thus if you run thousands of tests in a row, a lot of complex recomputation is done. This eventually becomes a performance issue. To speed this up - this project will calculate the entire $P_{N,k}$ for every value that is possible where $N < N_0$. Then the $P_{N,k}$ will be hardcoded inside the program for quick lookup of the value.

3.2 When can we approximate

The relative and absolute error behave very differently. The absolute error is calculated as follows.

$$\Delta = P_{k,N} - P_{approx,k,N} \quad (8)$$

The relative error is calculated as follows.

$$\epsilon_k = \frac{\Delta}{P_{approx,k,N}} \cdot 100\% \quad (9)$$

Note that the errors are signed, because this lets us know if the p-value estimation is over- or underestimated. The absolute error between accurate table values P and approximate table P_{approx} values first falls and then rises until at the center of the distribution, but gets stable and almost constant from the middle to the edge of the distribution. It can be seen on the Figure 5. Note that the error fluctuates. This is artifact of the V-table quickly fluctuating values that the smooth approximation cannot follow. The relative error gets smaller, the higher the number of samples N is. However, the relative error looks like a slide. It is stable at the center of the distribution, but starts rising sharply at the middle. At the edge of the distribution, the relative error is nearly 100%. It can be seen at the Figure 6.

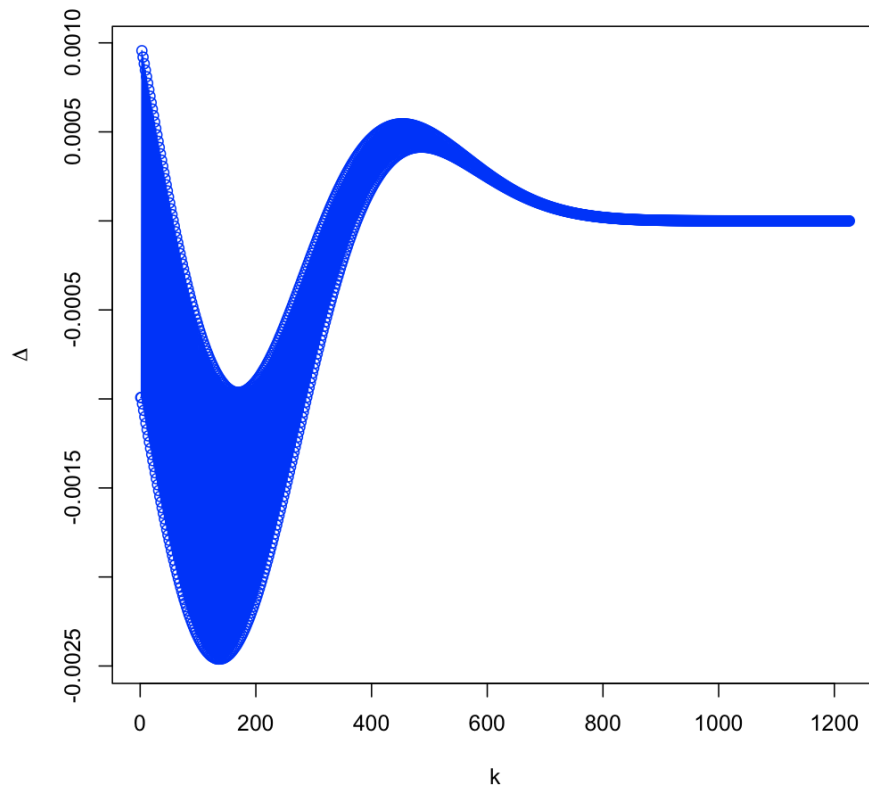


Figure 5: Absolute error between accurate P and approximate P_{approx} values for sample size $N = 50$.

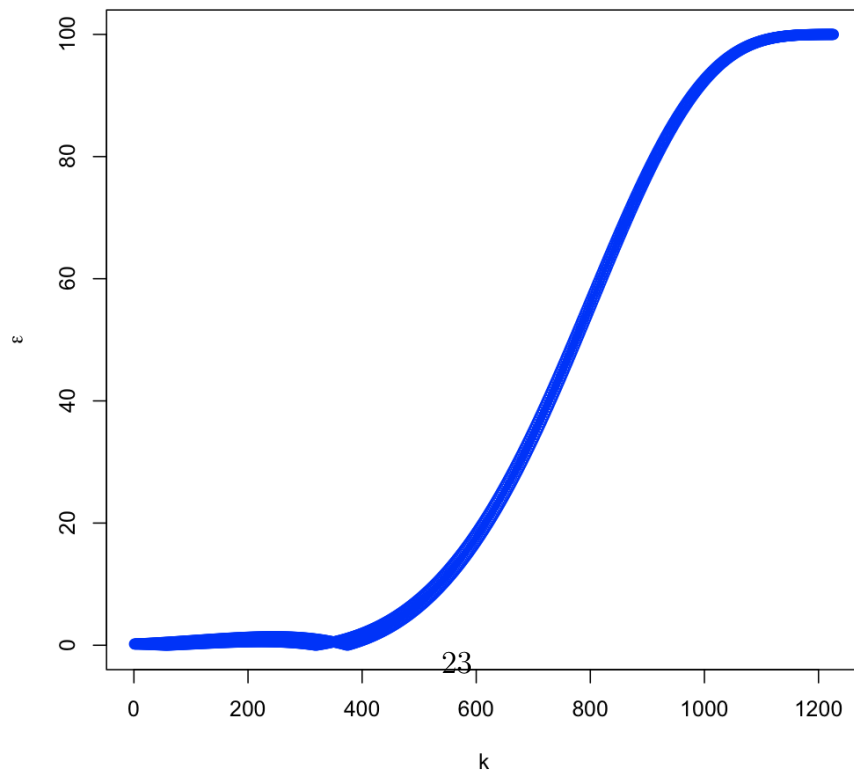


Figure 6: Relative error $|\epsilon|$ between accurate P and approximate P_{approx} values for sample size $N = 50$.

It is important to have low relative error on the region, i.e., good region, because this ensures that the test results stay accurate. If the relative error is high, then this will also increase the possibility of a false positive which is, rejecting H_0 when we actually should not.

First we wanted to confirm that as the number of samples grows, accurate table values P and approximate table values P_{approx} will get become more similar. To do this, we used the following equation for each N .

$$k_0 = \max_k \{k : P_{N,k} > 0\} \quad (10)$$

And compared them to the results as shown in the equation below, where ϵ_k is the relative error with the corresponding k value as shown in the equation (9).

$$k_1 = \max_k \{k : |\epsilon_k| \leq 5\%\} \quad (11)$$

The results of the table can be seen on Figure 7.

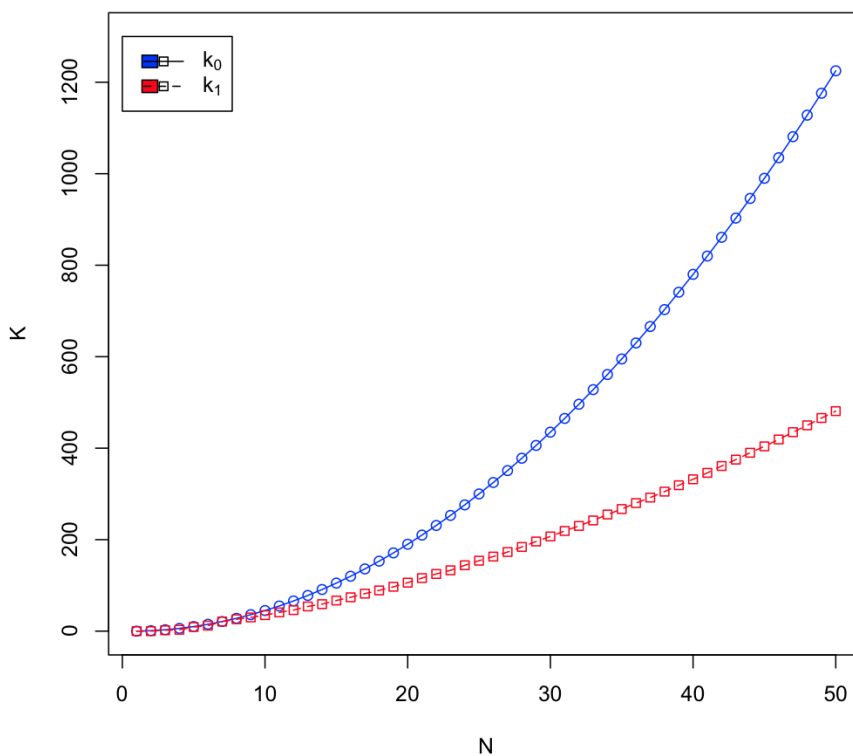


Figure 7: Region between low relative error and actual distribution edge increasing as sample size increases.

Much to our surprise, as N grew, the gap between approximate and accurate values grew bigger. This means that as N grows, the relative approximate error of Gaussian distribution will get more inaccurate toward the tails of the distribution. Previously we thought that with the growth of N , Gaussian would surely get more and more accurate. From Figure 8 we cannot really make a difference between accurate P and approximate P_{approx} values, but if we put the values to logarithmic scale, the difference is easier to notice. From Figure 9 we can see that approximated P_{approx} values are a little bit bigger toward the edge of distribution - as can be seen when k grows.

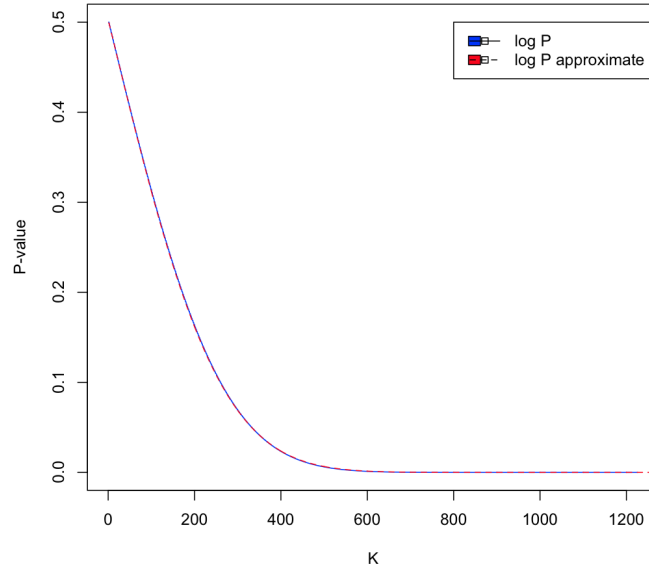


Figure 8: Accurate P and approximate P_{approx} values as k increases when sample size $N = 40$.

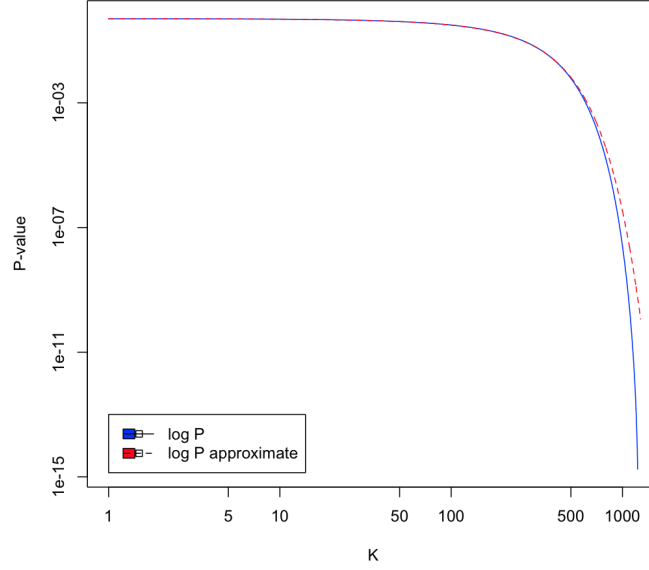


Figure 9: Logarithmic accurate P and approximate P_{approx} values as logarithmic k increases when sample size $N = 40$.

To further investigate this finding, we decided to find out the minimum approximate P value that we can get for each N while maintaining a certain relative error ϵ threshold. Let ϕ be the threshold. For that we first find the set.

$$\mathcal{K} = \{k : \epsilon_k < \phi\} \quad (12)$$

Then we choose the minimal approximate P value.

$$p = \min\{P_{approx,k} : k \in \mathcal{K}\} \quad (13)$$

The thresholds chosen were $\epsilon = 5\%$, $\epsilon = 10\%$, $\epsilon = 20\%$, $\epsilon = 50\%$. From Figure 10 we can see that as the N grows, the graph becomes stable, meaning that the distribution becomes a Gaussian distribution at around $N > 10$ and $N < 25$. When $N < 10$, then the approximation is so random that it cannot be used reliably. As an example of reading the graph, if $N = 30$, we can read from this graph that relative error $\epsilon < 5\%$ is attainable only for p-values that are larger than 0.012. We can also see from the figure that relative error ϵ does not increase monotonically when k grows, that's why there are fluctuations on the graph. We can see that the minimum P value under error threshold does get smaller, as N increases, so this further proves that Gaussian distribution does get more accurate as N grows. This is because as N increases, we can take a P value closer to the tail of the distribution and still get an accurate value.

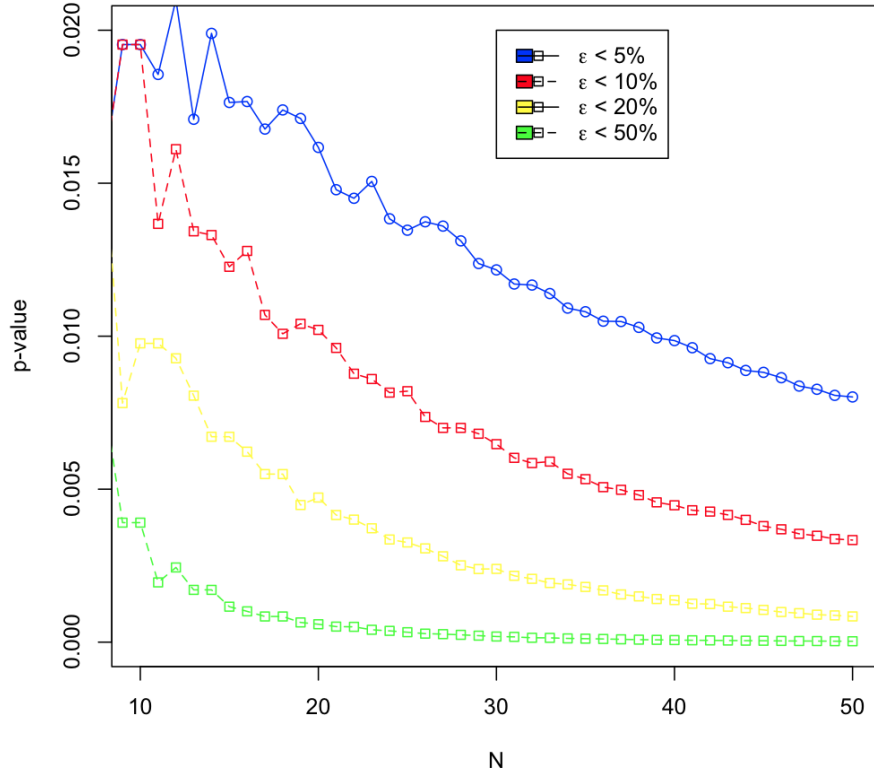


Figure 10: Relative error $|\epsilon|$ as sample size N increases

The Figure 11 shows us the relative error as P grows when $N = 20$. The P values are logged on this graph. There are two noteworthy things here, however. First, around $P = 0.05$, the relative error actually gets a little tip toward accuracy. Second, it can be seen that there are two accuracy paths that the error takes, depending on whether in $P_{approxN,K}$ the K is even or odd number. This is a computational artifact caused by the recurrence of the $V_{N,K}$ values and we will not make any conclusions of that. To understand this better, one must familiarize himself with the algorithm defined in the Chapter 2.7.2.

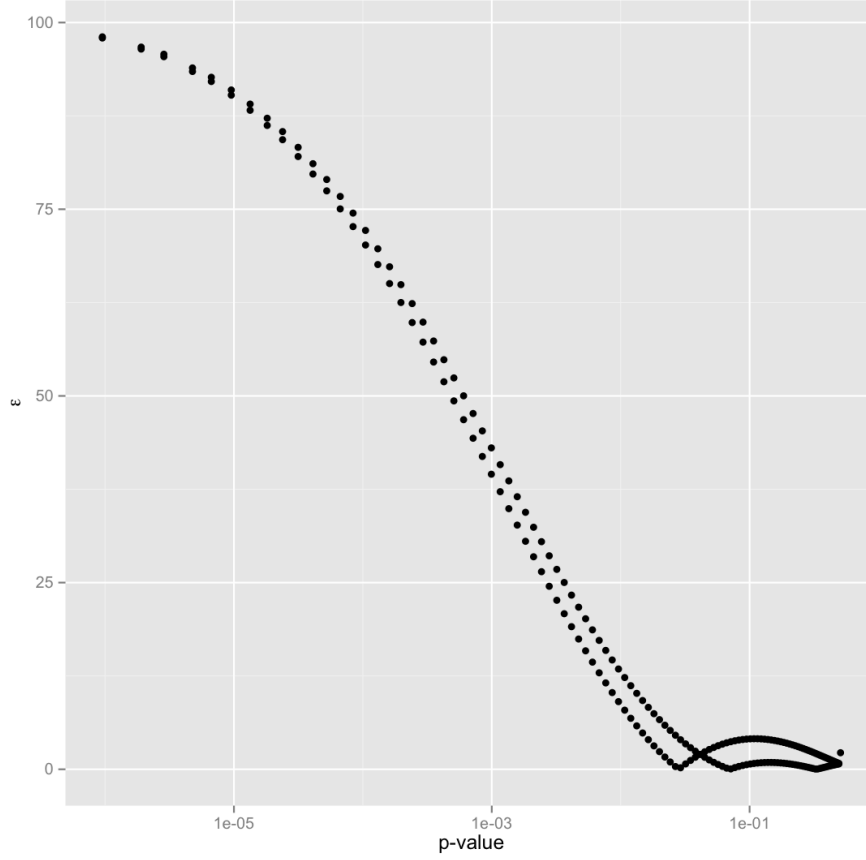


Figure 11: Relative error $|\epsilon|$ decreasing as p-value grows, when $N = 20$

3.3 The behavior of relative error

To further prove that we can start using Gaussian distribution at a certain N , we needed to prove that P_{approx} will get more accurate as the N increases toward the tail of the distribution as well. To achieve this, we found the largest relative error ϵ between P and P_{approx} for each N and for each p-value in the regions $[1, 0.1]$, $[0.1, 0.01]$, ..., $[10^{-6}; 10^{-7}]$. As can be seen from Figure 12 and Figure 13, the relative error increases fast. When p-value is between 10^{-5} and 10^{-6} the relative error is massive. When $N = 200$, then the error is still around 40%. The bigger the N gets, the bigger the error at the tail. However, it can also be seen that the relative error decreases steadily as N grows for each of the p-value ranges chosen. This means that even though the error of the distribution tail edge does increase as N grows, the approximate P_{approx} and accurate P values do get more similar as the N grows.

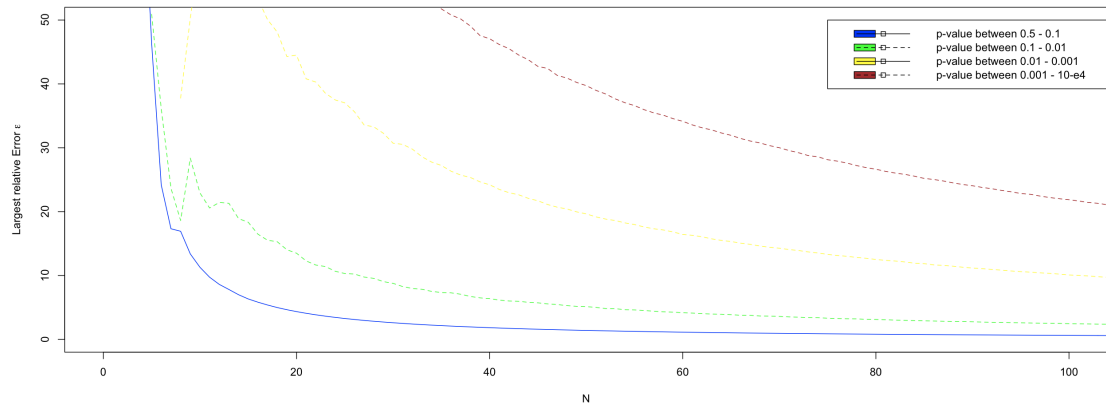


Figure 12: Showing the maximum relative error ϵ decreasing in probability ranges, as sample size N increases

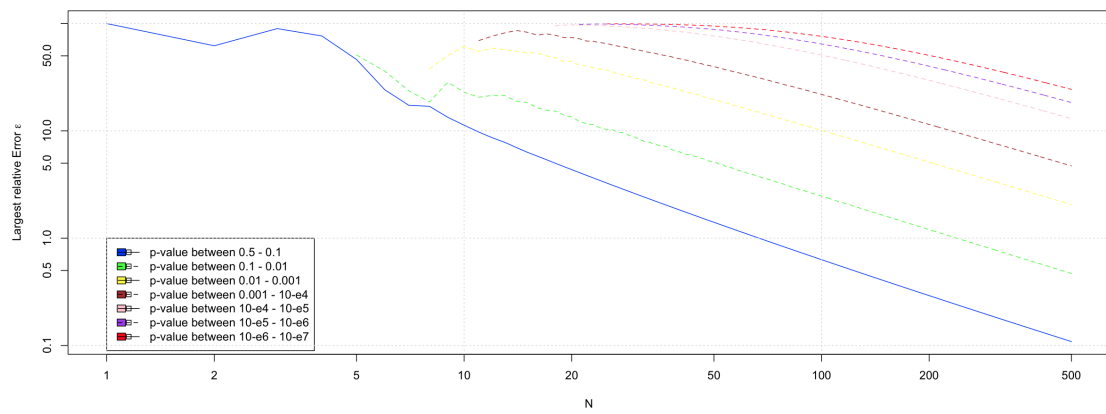


Figure 13: Showing the maximum relative error ϵ decreasing in probability ranges, as sample size N increases, when axes are log scaled

To investigate the relation between accurate P and approximate P_{approx} values, we wanted to see the difference between all p-values values when $N = 50$ and when $N = 25$. As can be seen from Figure 14 or Figure 15, the P_{approx} is constantly a little bit larger. When $N = 25$, the graph is a little less smooth than when $N = 50$. This follows the same story as before - the bigger the N , the more similar the p-values.

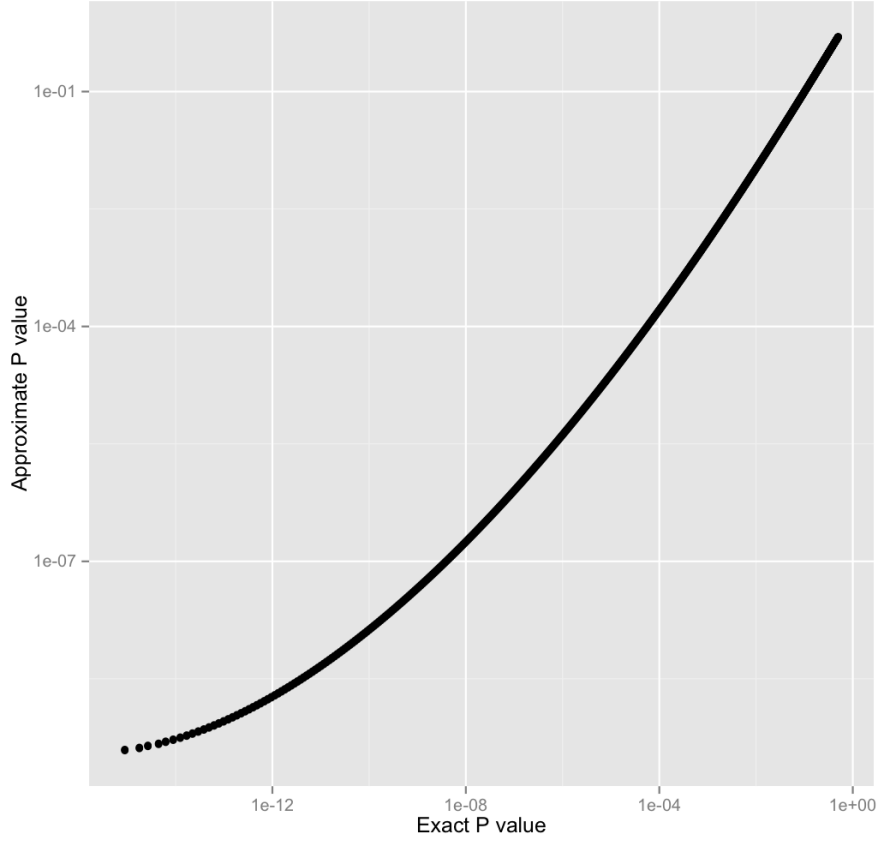


Figure 14: Relative value of all accurate P and approximated P_{approx} values when sample size N is 50

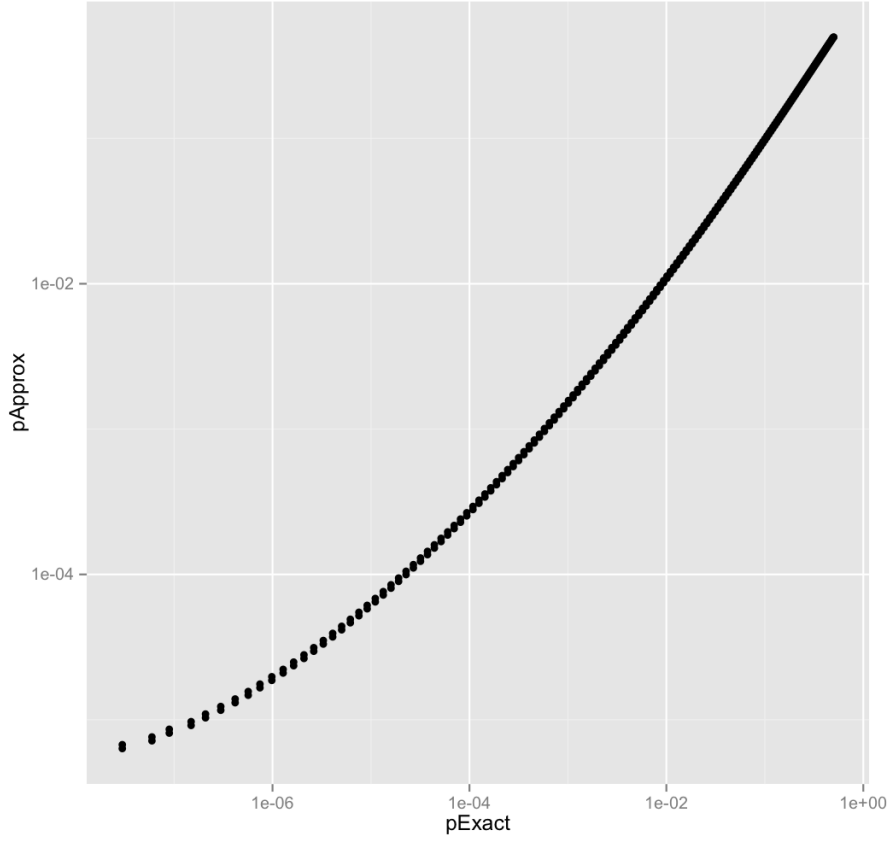


Figure 15: Relative value of all accurate P and approximated P_{approx} values when sample size N is 25

3.4 Finding out when can we approximate p-value

Judging by all the data that has been gathered, it can be clearly see that as N grows, accurate P and approximate P_{approx} values grow apart at the edges of distribution. however, overall the distribution gets more and more similar. The most helpful graph for determining the necessary number of samples N_0 we need to approximate p-value is Figure 13. Table 3 illustrates the number of samples N needed for a certain relative accuracy. So the number of samples N_0 we need to approximate depends on the number of measurements and the required accuracy.

Table 3: Showing the minimal required sample size for a certain relative error and number of measurements

| Measurements | Error | Required N |
|--------------|-------|------------|
| 10 | 5% | 25 |
| 100 | 5% | 250 |
| 1000 | 5% | 500 |
| 10000 | 5% | > 500 |
| 100000 | 5% | > 500 |
| 10 | 10% | 80 |
| 100 | 10% | 200 |
| 1000 | 10% | > 500 |
| 10000 | 10% | > 500 |
| 100000 | 10% | > 500 |
| 10 | 20% | 50 |
| 100 | 20% | 100 |
| 1000 | 20% | 280 |
| 10000 | 20% | 500 |
| 100000 | 20% | > 500 |
| 10 | 50% | 25 |
| 100 | 50% | 50 |
| 1000 | 50% | 100 |
| 10000 | 50% | 150 |
| 100000 | 50% | 200 |

4 Further optimizations

At the beginning of this paper, we raised the goal of optimizing our implementation on Wilcoxon test by hardcoding the accurate P value table $P_{N,K}$, so the program would not need to recompute it all the time. When given $N = 80$, then the table $P_{N,K}$ size was 2.5 MB. As pointed out in the beginning, when the N grows, the table $P_{N,k}$ size grows with it with the speed $O(n^2)$. At $N = 500$, the file size was over 5 GB. This raised the need to optimize the table $P_{N,k}$ further.

As the research found out, the required sample size N_0 to approximate p-value was very large and thus the corresponding p-value table is really large. To make it more practical, we need a way to pack the p-value table further. For that we used several methods. We noticed that the ration between true P values and approximated P_{approx} values (Equation (14)) is 1 until around the middle-range

of the P value table.

$$\frac{P_{k,N}}{P_{approx,k,N}} \quad (14)$$

Hence it makes sense to approximate equation (14) by some function γ_k and computer the approximation of a p-value as shown on equation (15).

$$P_{approx,k,N} \cdot \gamma_k \quad (15)$$

4.0.1 dose-response curve

Initially, when looking at the shape of ratio (14), it was clear that sigmoid functions used to represent dose-response curves should approximate the ratio pretty well. Since the curve is determined by up to 5 parameters, it was a good candidate for approximating the ration. Unfortunately, experiments showed that the data range, as shown in Figure 16, changes a lot toward the edge of the range. Further experiments showed that the algorithm was completely inaccurate at the far edge, so in conclusion it only worked for the middle range and not for the edges of the distribution [JsCr13].

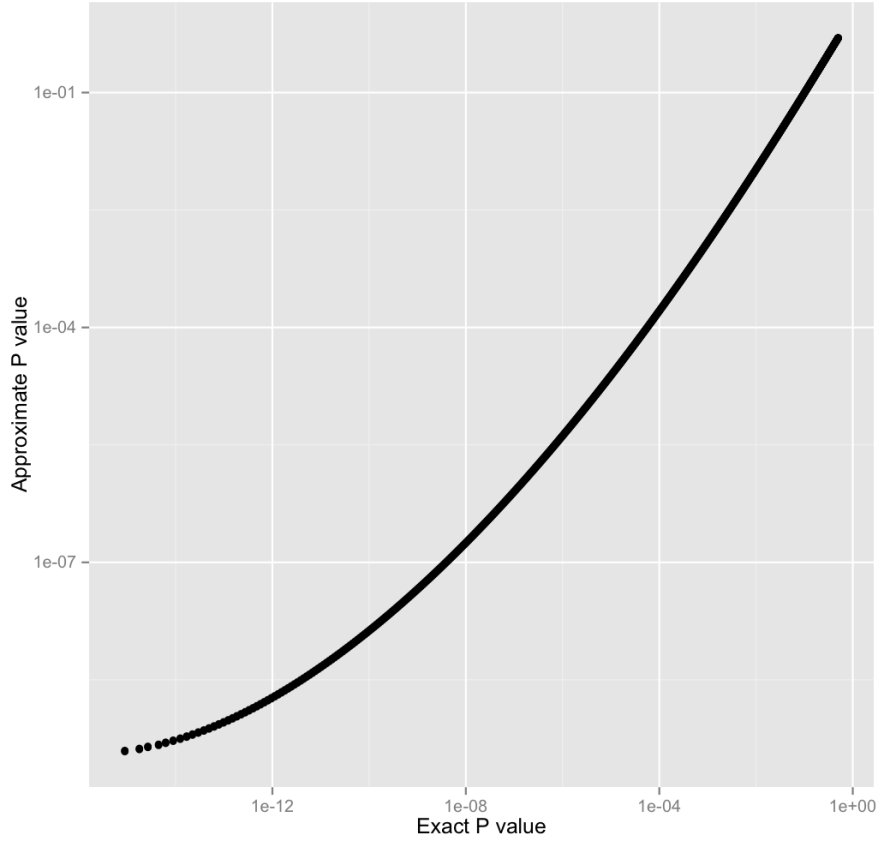


Figure 16: P range with DRC(sigmoid approximation) algorithm when $N=500$

4.1 Linear approximation

4.1.1 Goal

Since we did not want to do anything really clever, we used linear approximation as the second alternative. Our goal is to drop as many intermediate points from the P -table as we can, if they can be approximated with reasonable relative precision using linear interpolation.

4.1.2 Linear approximation

The linear approximation algorithm works by starting from a data point in the data range and then one by one moving to the next data point, until the points between them cannot be approximated with enough accuracy. Then the last accurate point is saved, picked as the next starting point and the process is repeated until the end of range has been reached. This guarantees approximated data range within

error margin. The graphical illustration of how linear approximation removes data points while staying within a certain error threshold and taking the general data array shape can be seen on the Figure 17.

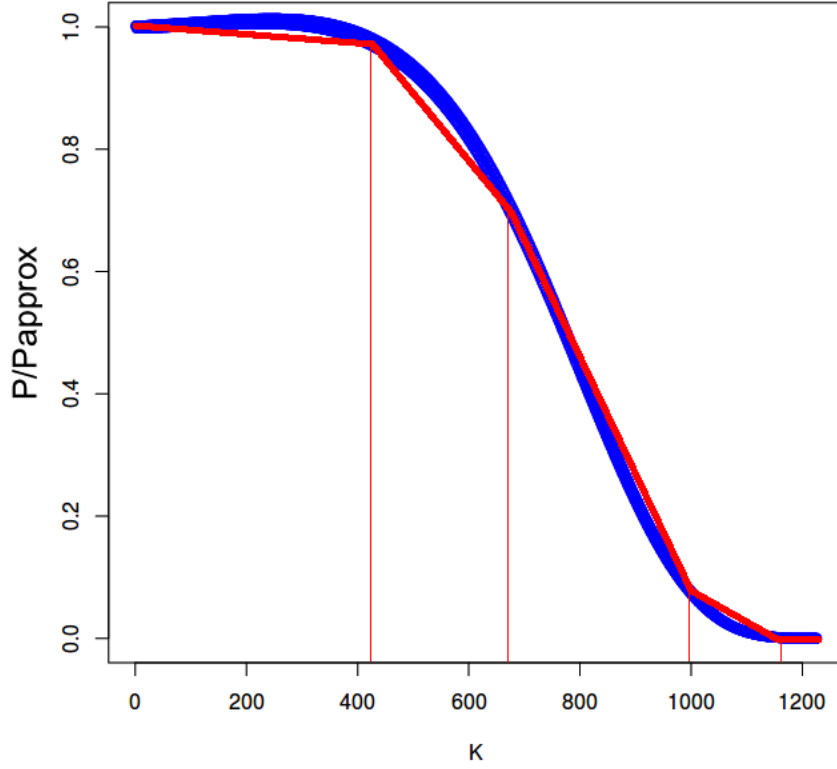


Figure 17: Linear approximation removing data points while staying within a certain error threshold and taking the general data array shape.

4.1.3 Linear interpolation

To reverse the algorithm, a very simple formula - linear interpolation is used. Let x be the approximated point k , x_1 the start points k , x_2 the end points k , y_1 the start points p-value, y_2 the end points p-value, y the approximated p-value [BLI2007].

$$y = y_1 + \frac{(y_2 - y_1) \cdot (x - x_1)}{(x_2 - x_1)} \quad (16)$$

A graphical explanation of the algorithm can be seen on the Figure 18.

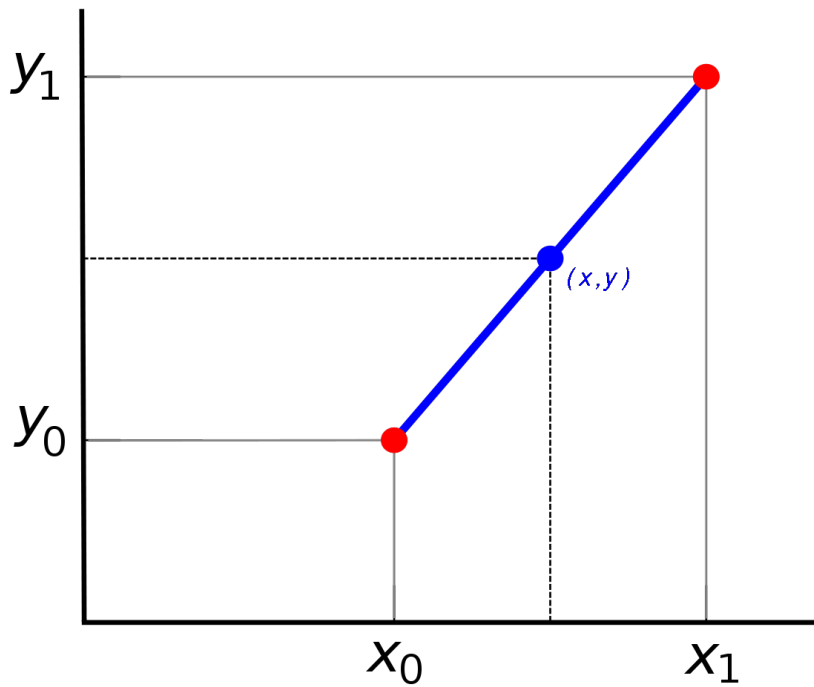


Figure 18: Linear interpolation finding y , when given x , start point and end point [BLI2007].

4.1.4 Algorithm for dropping intermediate points

The Linear approximation algorithm for dropping the data points in python looks as follows.

```
def linearInterpolate(approximatePoint,
                      endPoint,
                      startPoint,
                      endValue,
                      startValue):
    return startValue +
        ((endValue - startValue) * (approximatePoint -
                                     startPoint) /
         (endPoint - startPoint))

def canApproximate(actualPoints, currentPoint, nextPoint):
    startValue = actualPoints[currentPoint]
    endValue = actualPoints[nextPoint]

    for i in range(1, nextPoint - currentPoint):
        approximatePoint = currentPoint + i
        trueValue = actualPoints[approximatePoint]
```

```

        approximateValue = linearInterpolate(approximatePoint,
                                              currentPoint,
                                              nextPoint,
                                              startValue,
                                              endValue)

    if approximateValue == 0:
        return True

    relativeaccuracy = trueValue / approximateValue

    if relativeaccuracy < 0.9 or relativeaccuracy > 1.1:
        return False
    return True

def getNextApproximatedPoint(actualPoints, currentPoint):
    nextPoint = currentPoint + 1

    while(canApproximate(actualPoints, currentPoint, nextPoint)):
        if(nextPoint == len(actualPoints) - 1):
            break
        nextPoint += 1
    return nextPoint

def calculateApproximateRow(actualPoints):
    approximatedPoints = []
    currentPoint = 0
    approximatedPoints.append([currentPoint, actualPoints[
        currentPoint]])

    while(currentPoint != len(actualPoints) - 1):
        nextPoint = getNextApproximatedPoint(actualPoints,
                                              currentPoint)
        approximatedPoints.append([nextPoint, actualPoints[
            nextPoint]])
        currentPoint = nextPoint
    return approximatedPoints

```

4.2 Optimization summary

The tests showed that using relative P values between the P tables gave less approximated data points, than using straight forward accurate P table. In the end, thanks to linear approximation, the approximated table, when $N = 500$ was created. The error margin for the chosen approximated table was chosen to be 20%, since that was deemed as accurate enough. It is still a lot more accurate than Gaussian P table at the edges of the table. The approximated table is only 1.6 MB large. Smaller than the 2.5 MB accurate P table was when $N = 80$.

4.3 Further optimizations

A lot of further optimizations could be done on this subject, for example.

1. The algorithm to approximate accurate P table could be improved.
2. The plain text file could be compressed.
3. The shared library algorithm implementations could be improved and expanded.

5 The implementation

5.1 Library description

The implemented library consist of a C++ shared library, C++ terminal interface and R extended with C++ interface package.

The entire repository takes 8.1 MB, the shared library size is 1.9 MB, the terminal interface takes 1.9 MB, R interface takes 713 KB of disk space.

The performance tests show that the shared library could run over 20000 Wilcoxon tests with 120 samples in under 0.7 seconds. Compared to the vanilla R implementation which took many seconds to run 20 tests, this is a significant improvement.

5.2 Repository description

Repository that contains everything about our findings can be found in the git repository [JSrepo14]. The repository contains a number of folders.

WilcoxonTestLibrary. Implementation for the optimized Wilcoxon test. The library can be installed by going to the WilcoxonTestLibrary folder and running the command

```
|| $make install
```

RcppWilcoxonTest. An interface that connects our implementation of the optimized Wilcoxon algorithm to the R. Note that you must have installed WilcoxonTestLibrary to use it. The interface must be compiled with separate R commands from the command line. You need Rcpp packages for R to use it. It can be installed in R console by running [RFDE14]:

```
|| $install.packages("Rcpp")
```

After that, our Wilcoxon signed-rank implementation package can be installed to R by going to the RcppWilcoxonTest folder and running the command in the command line in the folder [ATSD14]:

```
|| $R CMD INSTALL .
```

You can now open the R console and load our library by running the command:

```
|| >library('RcppWilcoxonTest')
```

To invoke the Wilcoxon test function, run:

```
|| >RcppWilcoxonTest::WilxTest(dataMatrix, testIndexes,  
    controlIndexes)
```

TerminalWilcoxonTest. An interface that connects our implementation of the optimized Wilcoxon algorithm to the R. To use the terminal Wilcoxon test, you must first install Wilcoxon Test Library. It currently only supports NetCDF file as input data. The folder also contains example NetCDF file and the “help” gives instructions how to use it. The interface can be compiled by going to the TerminalWilcoxonTest folder and running the command:

```
|| $make
```

Further information on how to use it can be found out with the command

```
|| $/WilcoxonTest --help
```

WilcoxonVTable. Python program that can calculate V , P and approximated P tables, print them, create files of the tables and create a number of graphs on the tables. The python version used is 3. The main file is called main.py, which is in the WilcoxonVTable directory.

6 Conclusion

We found out that Gaussian distribution tail edge grows larger apart from the $P_{N,k}$ as N increases while the overall distribution grows more similar. In addition, we confirmed that if thousands of parallel tests are ran, then using approximation is not reliable and instead an accurate table of p-values must be used. Calculating that table is very expensive and thus we advise to pre-calculate the table for the program. However, the research shows that if 10 000 measurements are run under 5% relative error, then the hardcoded N should be over 1000. Table of that size would take many gigabytes of space and is impractical to use.

To get around that, we approximated the accurate P table. This allowed us to significantly reduce the size of the file and thus increase the hardcoded N that was implemented inside the library.

For the time being, we chose $N = 500$, with 20% error margin on the approximation. This is still significantly more accurate then Gaussian approximated table at the edges of the distribution. It will allow to make 1000 parallel measurement with 5% error margin.

The library itself could run over 20000 parallel tests with 100 samples in under a second. It is a shared C++ library that currently has implemented interfaces in GNU-R and Terminal.

Overall the research found some surprising results and can be considered a success.

References

- [FWil45] Wilcoxon, F. Individual Comparisons by Ranking Methods. International Biometric Society. *Biometrics Bulletin*, Vol. 1. Pages 80-83. 1945.
- [SSieg56] Siegel, S. Nonparametric statistics for the behavioral sciences. McGraw-Hill Humanities/Social Sciences/Languages; 2 edition. Pages 75-83. 1956.
- [RGSH11] Rew, R., Davis, G., Emmerson, S. and Daview, H. The NetCDF Users Guide. Unidata Program Center. Pages 5-7, 17-23. 2011.
- [RLow11] Lowry, R. Concepts & Applications of Inferential Statistics. Vassar College. Available from <http://vassarstats.net/textbook/ch12a.html>. ch.12a. 2011.
- [MTrio01] Triola, M. Elementary statistics (8 ed.). Addison Wesley. Pages 388. 2001.
- [JsCr13] Ritz, C. and Strebig, J. Dose-reponse curve. Available from <http://cran.r-project.org/web/packages/drc/index.html>. Pages 111. 2013.
- [ATSD14] DebRoy, S. and Trapletti, A. Writing R extensions. Free Software Foundation. Available from <http://cran.r-project.org/doc/manuals/R-exts.html>. ch 5-6. 2014.
- [RFDE14] Eddelbuettel, D. and Francois, R. Seamless R and C++ integration. Springer; 2013 edition. Available from <http://cran.r-project.org/web/packages/Rcpp/vignettes/Rcpp-introduction.pdf>. Pages 3-15. 2014.
- [RFDE14] Wikipedia Foundation Inc. linear interpolation. Available from http://en.wikipedia.org/wiki/Linear_interpolation. Last visited 10 January 2014.
- [RFDE14] Wikipedia Foundation Inc. Test Statistic. Available from http://en.wikipedia.org/wiki/Test_statistic. Last visited 10 January 2014.
- [ReWi14] Repapetillo, Wikipedia Foundation Inc. Illustration of p-value. Available from http://en.wikipedia.org/wiki/File:P-value_Graph.png. 24 March 2014.

- [JSpac14] Jerkku, S. Accurate P table calculation . Available from <https://github.com/stenver/wilcoxon-test/blob/master/WilcoxonVTable/pAccurateTable.py>. 2014.
- [JSrepo14] Jerkku, S. Project repository. Available from <https://github.com/stenver/wilcoxon-test/>. 2014.
- [BLI2007] Berland, Wikipedia Foundation Inc. Illustration of linear interpolation. Available from <http://en.wikipedia.org/wiki/File:LinearInterpolation.svg>. 2007.

Non-exclusive licence to reproduce thesis and make thesis public

I, Stenver Jerkku (date of birth: 10th of October 1990),

1. herewith grant the University of Tartu a free permit (non-exclusive licence) to:
 - 1.1 reproduce, for the purpose of preservation and making available to the public, including for addition to the DSpace digital archives until expiry of the term of validity of the copyright, and
 - 1.2 make available to the public via the web environment of the University of Tartu, including via the DSpace digital archives until expiry of the term of validity of the copyright,

Parallel Wilcoxon signed-rank tests

supervised by Sven Laur

2. I am aware of the fact that the author retains these rights.
3. I certify that granting the non-exclusive licence does not infringe the intellectual property rights or rights arising from the Personal Data Protection Act.

Tartu, 14.05.2014